

Universal Decoding Using a Noisy Codebook

Neri Merhav

The Andrew & Erna Viterbi Faculty of Electrical Engineering
Technion - Israel Institute of Technology
Technion City, Haifa 32000, ISRAEL
E-mail: merhav@ee.technion.ac.il

Abstract

We consider the topic of universal decoding with a decoder that does not have direct access to the codebook, but only to noisy versions of the various randomly generated codewords, a problem motivated by biometrical identification systems. Both the source that generates the original (clean) codewords, and the channel that corrupts them in generating the noisy codewords, as well as the main channel for communicating the messages, are all modeled by non-unifilar, finite-state systems (hidden Markov models). As in previous works on universal decoding, here too, the average error probability of our proposed universal decoder is shown to be as small as that of the optimal maximum likelihood (ML) decoder, up to a multiplicative factor that is a sub-exponential function of the block length. It therefore has the same error exponent, whenever the ML decoder has a positive error exponent. The universal decoding metric is based on Lempel-Ziv (LZ) incremental parsing of each noisy codeword jointly with the given channel output vector, but this metric is somewhat different from the one proposed in earlier works on universal decoding for finite-state channels, by Ziv (1985) and by Lapidoth and Ziv (1998). The reason for the difference is that here, unlike in those earlier works, the probability distribution that governs the (noisy) codewords is, in general, not uniform across its support. This non-uniformity of the codeword distribution also makes our derivation more challenging. Another reason for the more challenging analysis is the fact that the effective induced channel between the noisy codeword of the transmitted message and the main channel output is not a finite-state channel in general.

Index Terms: Universal decoding, finite-state channel, hidden Markov model, Lempel-Ziv algorithm, error exponent.

1 Introduction

The topic of universal decoding under channel uncertainty has received considerable attention in the last four decades. In [9] the *maximum mutual information* (MMI) decoder was first proposed and shown to achieve the capacity for discrete memoryless channels (DMC's). Csiszár and Körner [3] showed that the random coding error exponent of the MMI decoder, associated with a uniform random coding distribution over a given type class, achieves the same random coding error exponent as the maximum likelihood (ML) decoder. Csiszár [2] proved that for any modulo-additive DMC and the uniform random coding distribution over linear codes, the optimum random coding error exponent is universally achieved by a decoder that minimizes the empirical entropy of the difference between the output sequence and the input sequence. In [13], a parallel result was obtained for a certain class of memoryless Gaussian channels with slow fading and an unknown interference signal.

For channels with memory, Ziv [20] considered universal decoding for unknown unifilar finite-state (FS) channels with finite input and output alphabets, i.e., FS channels for which at each time instant, the next channel state is given by an unknown deterministic function of the channel current state, input and output. For ensembles of codes governed by the uniform distribution over a given permutation-invariant set of channel input vectors (namely, a type class or the disjoint union of several type classes), he proved that a decoder based on the Lempel–Ziv (LZ) incremental parsing algorithm asymptotically achieves the same error exponent as the ML decoder. In [11], Lapidoth and Ziv proved that the same universal decoder continues to be universally asymptotically optimum even for the broader class of FS channels with stochastic, rather than deterministic, next-state transitions. They still assumed a random coding distribution which is uniform over a given permutation-invariant set. In [7], Feder and Lapidoth have furnished sufficient conditions for general families of channels with memory to have universal decoders that asymptotically achieve the random coding error exponent of ML decoding. In [8], a competitive minimax criterion was proposed, in the quest for a more general systematic approach to the problem of universal decoding. Two additional related works on general methodologies for universal decoding are those of [12] and [14].

This paper is a further development on [11] and [20]. In particular, here we consider universal decoding in a situation where the decoder that does not have direct access to the codebook of the encoder, but only to noisy versions of the various randomly generated codewords, a problem motivated by applications in biometrical identification systems (see, e.g., [10, Section 5], [17], [18], [19], and many references therein) or other applications where storage, or finite-precision limitations do not enable the decoder to save the exact codewords of all messages, and then they must be quantized and hence distorted. In our model, both the source that generates the original (clean) codewords, and the channel that corrupts them in the process of generating the noisy codewords, as well as the main channel for communicating the messages, are all modeled by non-unifilar, FS systems (hidden Markov models). As in the previous above-mentioned works on universal decoding, here too, the average error probability of our proposed universal decoder is shown to be as small

as that of the optimal maximum likelihood (ML) decoder, up to a multiplicative factor that is a sub-exponential function of the block length, n . It therefore has the same error exponent, whenever the ML decoder has a positive error exponent. As in [11] and [20], the universal decoding metric is based on Lempel–Ziv (LZ) incremental parsing of each noisy codeword jointly with the given channel output vector, but this metric is somewhat different from that of [11] and [20]. Specifically, it includes an additional term, which is the logarithm of the induced probability of generating the noisy codeword of the message being tested. The reason for this difference is that here, unlike in [11] and [20], the probability distribution which governs the (noisy) codewords is, in general, not uniform across its support. This non-uniformity of the codeword distribution also makes our derivation quite more challenging. Another factor that makes the analysis here more involved is the fact that the effective induced channel between the noisy codeword of the transmitted message and the main channel output is not a FS channel in general.

The outline of the rest of the paper is as follows. In Section 2, we establish the notation conventions, define the problem formally, and spell out the assumptions. Section 3 is devoted to the statement of the main result and a discussion. Finally, in Section 4 the main results is proved.

2 Notation Conventions, Problem Formulation and Assumptions

2.1 Notation Conventions

Throughout the paper, random variables will be denoted by capital letters, specific values they may take will be denoted by the corresponding lower case letters, and their alphabets will be denoted by calligraphic letters. Random vectors and their realizations will be denoted, respectively, by capital letters and the corresponding lower case letters, both in the bold face font. Their alphabets will be superscripted by their dimensions. For example, the random vector $\mathbf{X} = (X_1, \dots, X_n)$, (n – positive integer) may take a specific vector value $\mathbf{x} = (x_1, \dots, x_n)$ in \mathcal{X}^n , the n -th order Cartesian power of \mathcal{X} , which is the alphabet of each component of this vector. The probability of an event \mathcal{E} (with respect to) w.r.t. a probability measure P will be denoted by $P[\mathcal{E}]$, and the expectation operator w.r.t. P will be denoted by $\mathbf{E}_P\{\cdot\}$. The subscript will be omitted if the underlying probability distribution is clear from the context. Logarithms and exponents will be defined w.r.t. the natural basis e , unless specified otherwise. In particular, $\exp_2(t)$ will sometimes be used to denote 2^t . The cardinality of a finite set, say, \mathcal{A} , will be denoted by $|\mathcal{A}|$.

2.2 Problem Formulation and Assumptions

Consider a coded communication system, defined as follows. First, a rate- R block code of length n , $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$, $M = e^{nR}$, is selected at random, where each $\mathbf{x}_m \in \mathcal{X}^n$, $m = 1, 2, \dots, M$, is drawn independently under a distribution $G(\mathbf{x})$. A message m is selected under the uniform distribution over the index set $\{1, 2, \dots, M\}$, and accordingly, the codeword \mathbf{x}_m is transmitted over a vector channel $W(\mathbf{z}|\mathbf{x})$, henceforth referred to as the *primary channel* (or the *main channel*), and the

resulting channel output vector, $\mathbf{z} \in \mathcal{Z}^n$, is received at the decoder side. The decoder, however, does not have access to the codebook, $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$, used by the encoder, but instead, it has access to a noisy version of that codebook, $\mathcal{C} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M\}$, $\mathbf{y}_m \in \mathcal{Y}^n$, $m = 1, 2, \dots, M$, where each \mathbf{y}_m is generated from the corresponding \mathbf{x}_m by another channel, $V(\mathbf{y}|\mathbf{x})$, henceforth referred to as the *secondary channel*. Clearly, this model, which was addressed by Willems *et al.* in [18] and [19] with application to biometrical identification systems (and later, further developed by Tuncel [17] and others), is formally equivalent to the ordinary model of channel random coding, where the codebook \mathcal{C} is selected at random, with each member, \mathbf{y}_m , being drawn independently under the random coding distribution,

$$P(\mathbf{y}) = \sum_{\mathbf{x} \in \mathcal{X}^n} G(\mathbf{x})V(\mathbf{y}|\mathbf{x}), \quad (1)$$

and where upon selecting the index m of the transmitted message, the corresponding codeword, \mathbf{y}_m , is transmitted over the channel

$$P(\mathbf{z}|\mathbf{y}) = \frac{P(\mathbf{y}, \mathbf{z})}{P(\mathbf{y})} = \frac{\sum_{\mathbf{x} \in \mathcal{X}^n} G(\mathbf{x})V(\mathbf{y}|\mathbf{x})W(\mathbf{z}|\mathbf{x})}{\sum_{\mathbf{x} \in \mathcal{X}^n} G(\mathbf{x})V(\mathbf{y}|\mathbf{x})}. \quad (2)$$

From this point onward, the original codebook $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$ no longer plays a role. Accordingly, we henceforth refer to $\{P(\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n\}$ as the *induced random coding distribution* (or the *effective random coding distribution*), and to $\{P(\mathbf{z}|\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n, \mathbf{z} \in \mathcal{Z}^n\}$ – as the *induced channel* (or the *effective channel*). Clearly, if G is a discrete memoryless source (DMS) and V is a discrete memoryless channel (DMC), then $\{P(\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n\}$ is a DMS as well. If, in addition, W is also a DMC, then so is the channel $\{P(\mathbf{z}|\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n, \mathbf{z} \in \mathcal{Z}^n\}$. In this case, the capacity of the system is simply the mutual information, $I(Y; Z)$, pertaining to the single-letter marginal $\{P(y, z), y \in \mathcal{Y}, z \in \mathcal{Z}\}$, see [18], [19]. It should be noted, however, that unlike the traditional model of random coding for channels, where random coding is a technical concept that merely serves the purpose of proving existence of good codes, here, when it comes to biometrical systems applications, the randomness of the code is part of the model setting. As a consequence, both G and V , and hence also the induced random coding distribution, $\{P(\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n\}$, are dictated to us, and are not subjected to our control.¹

As in [18], [19], here too, it is assumed that all three alphabets, \mathcal{X} , \mathcal{Y} , and \mathcal{Z} , are finite. In this paper, however, we go considerably beyond the realm of memoryless systems, and allow G , V and W to be all non-unifilar, FS systems (hidden Markov models), as follows. The distribution G assumes the form

$$G(\mathbf{x}) = \sum_{\boldsymbol{\omega}} \prod_{i=1}^n G(x_i, \omega_i | \omega_{i-1}), \quad (3)$$

where \mathbf{x} is as before, $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)$ is the source state vector, whose components take on values in a finite set Ω , and the initial state, ω_0 is assumed fixed. The primary channel, W , is modeled as

$$W(\mathbf{z}|\mathbf{x}) = \sum_{\boldsymbol{\sigma}} \prod_{i=1}^n W(z_i, \sigma_i | x_i, \sigma_{i-1}), \quad (4)$$

¹For this reason, the capacity is simply given by $I(Y; Z)$, without maximizing over the distribution of Y .

where $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_n)$ is the channel state vector, whose components take on values in a finite set Σ and the initial state, σ_0 , is fixed. Likewise, the secondary channel, V , is given by

$$V(\mathbf{y}|\mathbf{x}) = \sum_{\boldsymbol{\theta}} \prod_{i=1}^n V(y_i, \theta_i | x_i, \theta_{i-1}), \quad (5)$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ is the state vector whose components take on values in a finite set Θ and there is fixed initial state, θ_0 .

We consider the problem of universal decoding for the effective channel $P(\mathbf{z}|\mathbf{y})$ induced by the source (3), the main channel (4) and the secondary channel (5), according to (2). We will assume that G , V and W are not known to the decoder, and hence nor is the effective channel $\{P(\mathbf{z}|\mathbf{y}) \mid \mathbf{y} \in \mathcal{Y}^n, \mathbf{z} \in \mathcal{Z}^n\}$. Nonetheless, the effective random coding distribution, $\{P(\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n\}$, will be assumed known to the decoder. The rationale behind the latter assumption stems from the fact that the decoder knows the codebook, $\mathcal{C} = \{\mathbf{y}_1, \dots, \mathbf{y}_M\}$, and so, it has access to an exponential amount of data from which the parameters of this distribution can be estimated very accurately. In particular, note that $P(\mathbf{y})$ has a hidden Markov structure,

$$\begin{aligned} P(\mathbf{y}) &= \sum_{\mathbf{x}} G(\mathbf{x}) V(\mathbf{y}|\mathbf{x}) \\ &= \sum_{\boldsymbol{\theta}, \boldsymbol{\omega}, \mathbf{x}} \prod_{i=1}^n G(x_i, \omega_i | \omega_{i-1}) V(y_i, \theta_i | x_i, \theta_{i-1}) \\ &= \sum_{\boldsymbol{\theta}, \boldsymbol{\omega}} \prod_{i=1}^n \left[\sum_{\mathbf{x}} G(\mathbf{x}, \boldsymbol{\omega} | \boldsymbol{\omega}_{i-1}) V(y_i, \theta_i | \mathbf{x}, \theta_{i-1}) \right] \\ &= \sum_{\boldsymbol{\theta}, \boldsymbol{\omega}} \prod_{i=1}^n \pi(y_i, \theta_i, \omega_i | \theta_{i-1}, \omega_{i-1}), \end{aligned} \quad (6)$$

where in the last passage, we have defined the parameters $\pi(y, \theta, \omega | \theta', \omega') \triangleq \sum_{\mathbf{x}} G(\mathbf{x}, \boldsymbol{\omega} | \boldsymbol{\omega}') V(y, \theta | \mathbf{x}, \theta')$. These parameters can be estimated using well known estimation methods for hidden Markov models.² It will be assumed³ that

$$\pi(y, \theta, \omega | \theta', \omega') > 0 \quad (7)$$

for all $(\omega, \omega', \theta, \theta', y) \in \Omega^2 \times \Theta^2 \times \mathcal{Y}$, and we denote $\pi_{\min} \triangleq \min_{\omega, \omega', \theta, \theta', y} \pi(y, \theta, \omega | \theta', \omega')$.

Like in previous works on universal decoding, our objective is to devise a universal decoding metric whose average error probability is of the same exponential order as that of the ML decoder. As described in the Introduction, the problem of universal decoding for FS channels was considered

² The ML estimator for the parameters of a hidden Markov model, is known to be strongly consistent [1], [15]. More practically, one may use the iterative Baum algorithm, which is an instance of the EM algorithm [5] (see also the tutorials [6], [16] and references therein).

³ Note that this assumption concerns G and V only, it has nothing to do with the primary channel W . If the parameters $\{\pi(y, \theta, \omega | \theta', \omega')\}$ are estimated using the ML estimator (referring to footnote 2), then eq. (7) can be imposed as a constraint on the estimator.

first in [20], where it was assumed that the next-state transitions are given by a deterministic function of the current state, the current input and the current output. In [11], the framework was extended to handle general FS channels, where the state transitions were allowed to be stochastic (as in eqs. (4) and (5) above). Also, in both [11] and [20], the random coding distribution was assumed uniform across a given permutation-invariant set.⁴ Here the situation is different from both [11] and [20] because of two reasons.

1. The effective random coding distribution $\{P(\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n\}$ is not uniform over a permutation-invariant set, in general.
2. The effective channel $\{P(\mathbf{z}|\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n, \mathbf{z} \in \mathcal{Z}^n\}$ is not a FS channel, in general.

These differences are important, because in [11] and [20], both assumptions were used rather heavily.

For a given noisy code \mathcal{C} and a given channel output vector \mathbf{z} , let us define (similarly as in [7] and [11]) the ranking of the members of \mathcal{Y}^n , according to descending likelihood values, i.e., $P(\mathbf{z}|\mathbf{y}[1]) \geq P(\mathbf{z}|\mathbf{y}[2]) \geq \dots$, and let us denote by $M_o(\mathbf{y}, \mathbf{z})$ the ranking of \mathbf{y} given \mathbf{z} . For a given \mathbf{z} , the ranking function $M_o(\mathbf{y}, \mathbf{z})$ is therefore a one-to-one mapping from \mathcal{Y}^n to the set $\{1, 2, \dots, |\mathcal{Y}|^n\}$ with the property that $P(\mathbf{z}|\mathbf{y}') > P(\mathbf{z}|\mathbf{y})$ implies $M_o(\mathbf{y}', \mathbf{z}) < M_o(\mathbf{y}, \mathbf{z})$. The probability of error associated with the ML decoder for the given code \mathcal{C} and the effective channel, $\{P(\mathbf{z}|\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n, \mathbf{z} \in \mathcal{Z}^n\}$, is given by

$$P_{e,o}(\mathcal{C}) = \frac{1}{M} \sum_{m=1}^M P \left[\bigcup_{m' \neq m} \{M_o(\mathbf{y}_{m'}, \mathbf{Z}) \leq M_o(\mathbf{y}_m, \mathbf{Z})\} \middle| \text{message } m \text{ was sent} \right], \quad (8)$$

where the event $M_o(\mathbf{y}_{m'}, \mathbf{Z}) = M_o(\mathbf{y}_m, \mathbf{Z})$ accounts for the case where $\mathbf{y}_{m'} = \mathbf{y}_m$ (which is possible since the members of \mathcal{C} are chosen independently at random). The average probability of error w.r.t. the randomness of \mathcal{C} , is then

$$\bar{P}_{e,o} = \mathbf{E} \{P_{e,o}(\mathcal{C})\} \quad (9)$$

$$= 1 - \sum_{\mathbf{y}, \mathbf{z}} P(\mathbf{y}, \mathbf{z}) (1 - P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})])^{e^{nR}-1}, \quad (10)$$

where

$$\mathcal{E}_o(\mathbf{y}, \mathbf{z}) \triangleq \{\mathbf{y}' : M_o(\mathbf{y}', \mathbf{z}) \leq M_o(\mathbf{y}, \mathbf{z})\}. \quad (11)$$

As in [7] and [11], for later use, we define the function

$$f(t) \triangleq 1 - (1 - t)^{e^{nR}-1}, \quad t \in [0, 1], \quad (12)$$

and so,

$$\bar{P}_{e,o} = \sum_{\mathbf{y}, \mathbf{z}} P(\mathbf{y}, \mathbf{z}) f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]). \quad (13)$$

⁴A permutation-invariant set is a set that is closed under permutations, in other words, a set that can be represented by the disjoint union of type classes.

By the same token, for an arbitrary decoding metric $u(\mathbf{y}, \mathbf{z})$, we define a ranking function $M_u(\mathbf{y}, \mathbf{z})$, as any one-to-one mapping $\mathcal{Y}^n \rightarrow \{1, 2, \dots, |\mathcal{Y}|^n\}$ given \mathbf{z} , such that $u(\mathbf{y}', \mathbf{z}) < u(\mathbf{y}, \mathbf{z})$ implies $M_u(\mathbf{y}', \mathbf{z}) < M_u(\mathbf{y}, \mathbf{z})$. Accordingly, the average error probability associated with $u(\cdot, \cdot)$, is given by

$$\bar{P}_{e,u} = \sum_{\mathbf{y}, \mathbf{z}} P(\mathbf{y}, \mathbf{z}) f(P[\mathcal{E}_u(\mathbf{y}, \mathbf{z})]), \quad (14)$$

where

$$\mathcal{E}_u(\mathbf{y}, \mathbf{z}) \triangleq \{\mathbf{y}' : M_u(\mathbf{y}', \mathbf{z}) \leq M_u(\mathbf{y}, \mathbf{z})\}. \quad (15)$$

We are interested in a universal metric $u(\cdot, \cdot)$, that is independent of the unknown effective channel (but possibly dependent on the effective random coding distribution), such that $\bar{P}_{e,u}$ would not exceed $\bar{P}_{e,o}$ by more than a sub-exponential function of n , i.e.,

$$\bar{P}_{e,u} \leq e^{n\epsilon(n)} \bar{P}_{e,o}, \quad (16)$$

where $\epsilon(n) \rightarrow 0$ as $n \rightarrow \infty$.

3 Main Result

Given two sequences, \mathbf{y} and \mathbf{z} , both of length n , consider the joint incremental parsing [21] of the sequence of pairs

$$(y_1, z_1), (y_2, z_2), \dots, (y_n, z_n)$$

into c distinct phrases. Specifically, denoting $w_i = (y_i, z_i)$, $i = 1, 2, \dots, n$, we parse $\mathbf{w} = (w_1, \dots, w_n)$, sequentially into the distinct⁵ phrases, $w_1^{n_1}, w_{n_1+1}^{n_2}, \dots, w_{n_{c-1}+1}^n$, where $n_i + 1$ is the starting point of the i -th phrase, $i = 1, 2, \dots, c$ ($n_0 = 0$). According to the incremental parsing procedure of the LZ algorithm, each phrase $w_{n_i+1}^{n_{i+1}}$ is the shortest string that has not been encountered before as a parsed phrase, which means that its prefix, $w_{n_i+1}^{n_{i+1}-1}$, is identical to an earlier phrase, $w_{n_j+1}^{n_{j+1}}$, $j < i$. Let $c \equiv c(\mathbf{y}, \mathbf{z})$ denote the number of distinct phrases. For example,⁶ if

$$\begin{aligned} \mathbf{y} &= 0 \mid 1 \mid 0 \mid 0 \mid 0 \mid 1 \mid \\ \mathbf{z} &= 0 \mid 1 \mid 0 \mid 1 \mid 0 \mid 1 \mid \end{aligned}$$

then $c(\mathbf{y}, \mathbf{z}) = 4$. Let $c(\mathbf{z})$ denote the resulting number of distinct phrases of \mathbf{z} , and let $\mathbf{z}(\ell)$ denote the ℓ th distinct \mathbf{z} -phrase, $\ell = 1, 2, \dots, c(\mathbf{z})$. In the above example, $c(\mathbf{z}) = 3$. Denote by $c_\ell(\mathbf{y}|\mathbf{z})$ the number of occurrences of $\mathbf{z}(\ell)$ in the parsing of \mathbf{z} , or equivalently, the number of distinct \mathbf{y} -phrases that jointly appear with $\mathbf{z}(\ell)$. Clearly, $\sum_{\ell=1}^{c(\mathbf{z})} c_\ell(\mathbf{y}|\mathbf{z}) = c(\mathbf{y}, \mathbf{z})$. In the above example, $\mathbf{z}(1) = 0$, $\mathbf{z}(2) = 1$, $\mathbf{z}(3) = 01$, $c_1(\mathbf{y}|\mathbf{z}) = c_2(\mathbf{y}|\mathbf{z}) = 1$, and $c_3(\mathbf{y}|\mathbf{z}) = 2$. We next define our universal decoding metric as

$$u(\mathbf{y}, \mathbf{z}) \triangleq \log P(\mathbf{y}) + \sum_{\ell=1}^{c(\mathbf{z})} c_\ell(\mathbf{y}|\mathbf{z}) \log c_\ell(\mathbf{y}|\mathbf{z}), \quad (17)$$

⁵ To be more precise, the phrases are all distinct with the possible exception of the last phrase, which might be incomplete.

⁶The same example appears in [20].

which in turn, defines the decoder

$$\hat{m}_u = \arg \min_m u(\mathbf{y}_m, \mathbf{z}), \quad (18)$$

where ties broken according to an arbitrary ranking function $M_u(\cdot, \mathbf{z})$ associated with (17).

We are now ready to state our main result, whose proof appears in Section 4.

Theorem 1 *Under the assumptions of Subsection 2.2, the universal decoder (18) satisfies eq. (16) where $\epsilon(n) = O((\log \log n)/\log n)$, with a leading term⁷ that is linear in $\log |\mathcal{Y} \times \mathcal{Z}|$.*

It should be noticed that the universal decoding metric (17) is different from the one in [11] and [20], because it includes the term $\log P(\mathbf{y})$ in addition to the LZ conditional compressibility term, $\sum_{\ell=1}^{c(\mathbf{z})} c_\ell(\mathbf{y}|\mathbf{z}) \log c_\ell(\mathbf{y}|\mathbf{z})$ (see also [14]). The reason for this difference is that the effective random coding distribution, $\{P(\mathbf{y}), \mathbf{y} \in \mathcal{Y}^n\}$, is not necessarily uniform over its support, in contrast to the assumption in both [11] and [20]. In a way, the decoder (18) can be seen as an extension of the MMI decoder, which is the well known universal decoder for DMCs [3]. To see this, observe that (18) can be rewritten as

$$\hat{m}_u = \arg \max_m \left\{ \frac{1}{n} \log \left[\frac{1}{P(\mathbf{y}_m)} \right] - \frac{1}{n} \sum_{\ell=1}^{c(\mathbf{z})} c_\ell(\mathbf{y}_m|\mathbf{z}) \log c_\ell(\mathbf{y}_m|\mathbf{z}) \right\}, \quad (19)$$

where the term $\frac{1}{n} \log[1/P(\mathbf{y}_m)]$ plays a role like the empirical entropy associated with \mathbf{y}_m and the term $\frac{1}{n} \sum_{\ell=1}^{c(\mathbf{z})} c_\ell(\mathbf{y}_m|\mathbf{z}) \log c_\ell(\mathbf{y}_m|\mathbf{z})$ is parallel to the conditional empirical entropy of \mathbf{y}_m given \mathbf{z} . Thus, the difference is analogous to a certain notion of a generalized empirical mutual information. But having said that, we should add a digression that, when confining the discussion to the memoryless case, the first term in (19) gives the empirical entropy of \mathbf{y}_m only in the case where $\{P(\mathbf{y})\}$ is uniform across a single type class. If instead, it is a product distribution, then the MMI metric should be supplemented with a divergence term between the empirical distribution and the true distribution.⁸

The proof of Theorem 1 contains essentially similar ingredients to those in [11]. There are, however, a few differences that should be pointed out. In the previous paragraph, we mentioned that here, as opposed to those papers, the random coding distribution is not uniform in general. This difference is also responsible for the fact that there are a few non-trivial issues in the extension of the derivations of [11] and [20] to our setting, as in those two earlier papers, the uniformity of the random coding distribution (across its support), was used quite heavily. In particular, the pairwise error probability, $P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]$, which plays a central role in the analysis in [11] and [20], is simply proportional to the cardinality of $\mathcal{E}_o(\mathbf{y}, \mathbf{z})$, namely to $M_o(\mathbf{y}, \mathbf{z})$, which in turn, can be evaluated using combinatorial considerations. Here, on the other hand, the members of $\mathcal{E}_o(\mathbf{y}, \mathbf{z})$ have to

⁷The sequence $\epsilon(n)$ depends also on other parameters of the problem, like $|\Theta|$, $|\Omega|$, $|\Sigma|$, and π_{\min} , but these parameters appear in negligible terms of $\epsilon(n)$, that decay faster than $(\log \log n)/\log n$.

⁸In this context, the author has some doubts concerning the asymptotic optimality of the MMI decoder used in [4].

be weighed by their various probabilities, $\{P(\mathbf{y}'), \mathbf{y}' \in \mathcal{E}_o(\mathbf{y}, \mathbf{z})\}$. In particular, in an important technical lemma of [11] (Lemma 2 therein), the last step of the proof is relatively easy, because thanks to the uniformity assumption therein, it is associated with the calculation of the quantity, $\sum_{\mathbf{y}} 1/M_o(\mathbf{y}, \mathbf{z})$ (in our notation), which is nothing but the harmonic series, $\sum_{i=1}^N 1/i \leq \ln N + 1$ (N – positive integer), as $M_o(\mathbf{y}, \mathbf{z})$ is defined as a ranking function (see, in particular, the last step in the chain of inequalities at the end of page 1751 in [11]). For the non-uniform input considered here, the relevant extension of the above mentioned expression turns out to be $\sum_{\mathbf{y}} P(\mathbf{y})/P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]$, which is not as straightforward to bound in a useful manner. Fortunately enough, as is shown in Lemma 1 below, this can nevertheless still be done, and in a quite general manner, that is almost completely unrelated to the hidden Markov structure of the model. Another source for some technical challenges is the fact that the induced channel, $\{P(\mathbf{z}|\mathbf{y})\}$, is not a FS channel, in general. This calls for separate treatment of the numerator and the denominator of $P(\mathbf{z}|\mathbf{y}) = P(\mathbf{y}, \mathbf{z})/P(\mathbf{y})$ (which both obey a hidden Markov model), that in turn, may be dominated by two different sequences of states. Nonetheless, these difficulties can also be circumvented, as will be seen in Section 4.

4 Proof of Theorem 1

The idea of the proof is to lower bound $\bar{P}_{e,o}$ and to upper bound $\bar{P}_{e,u}$ by two expressions which are identical up to a multiplicative factor of $e^{n\epsilon(n)}$. We begin with the upper bound to $\bar{P}_{e,u}$.

Let us denote

$$v(\mathbf{y}, \mathbf{z}) \triangleq \sum_{\ell=1}^{c(\mathbf{z})} c_{\ell}(\mathbf{y}|\mathbf{z}) \log c_{\ell}(\mathbf{y}|\mathbf{z}), \quad (20)$$

so that $u(\mathbf{y}, \mathbf{z}) = \log P(\mathbf{y}) + v(\mathbf{y}, \mathbf{z})$. We will use the fact that $v(\mathbf{y}, \mathbf{z})$ is almost large enough to serve as a legitimate length function for lossless compression of \mathbf{y} given \mathbf{z} , where \mathbf{z} serves as side information available to both the encoder and the decoder. In particular, in the proof of Lemma 2 in [20, p. 460], Ziv describes a lossless compression scheme with side information, whose length function, $L(\mathbf{y}|\mathbf{z})$, satisfies

$$L(\mathbf{y}|\mathbf{z}) \leq v(\mathbf{y}, \mathbf{z}) + n\epsilon_1(n), \quad (21)$$

with

$$\epsilon_1(n) = O\left(\frac{\log \log n}{\log n}\right), \quad (22)$$

whose leading term is linear in $\log |\mathcal{Y} \times \mathcal{Z}|$. Now, let us define

$$\bar{P}_{e,u}(\mathbf{y}, \mathbf{z}) = f(P[\mathcal{E}_u(\mathbf{y}, \mathbf{z})]), \quad (23)$$

where $f(\cdot)$ is defined as in (12). Now,

$$P[\mathcal{E}_u(\mathbf{y}, \mathbf{z})] = \sum_{\{\mathbf{y}': M_u(\mathbf{y}', \mathbf{z}) \leq M_u(\mathbf{y}, \mathbf{z})\}} P(\mathbf{y}')$$

$$\begin{aligned}
&\leq \sum_{\{\mathbf{y}': P(\mathbf{y}') \exp_2[v(\mathbf{y}', \mathbf{z})] \leq P(\mathbf{y}) \exp_2[v(\mathbf{y}, \mathbf{z})]\}} P(\mathbf{y}') \\
&\leq \sum_{\{\mathbf{y}': P(\mathbf{y}') \exp_2[v(\mathbf{y}', \mathbf{z})] \leq P(\mathbf{y}) \exp_2[v(\mathbf{y}, \mathbf{z})]\}} P(\mathbf{y}) \exp_2[v(\mathbf{y}, \mathbf{z}) - v(\mathbf{y}', \mathbf{z})] \\
&\leq P(\mathbf{y}) \exp_2[v(\mathbf{y}, \mathbf{z})] \sum_{\mathbf{y}' \in \mathcal{Y}^n} \exp_2[-v(\mathbf{y}', \mathbf{z})] \\
&\leq 2^{n\epsilon_1(n)} P(\mathbf{y}) \exp_2[v(\mathbf{y}, \mathbf{z})] \sum_{\mathbf{y}' \in \mathcal{Y}^n} 2^{-L(\mathbf{y}'|\mathbf{z})} \\
&\leq e^{n\epsilon_1(n)} P(\mathbf{y}) \cdot \exp_2[v(\mathbf{y}, \mathbf{z})] \\
&= e^{n\epsilon_1(n)} \cdot \exp_2[u(\mathbf{y}, \mathbf{z})], \tag{24}
\end{aligned}$$

where in the second to the last step, we have used Kraft's inequality and we bounded $2^{n\epsilon_1(n)}$ by $e^{n\epsilon_1(n)}$, simply for convenience in later steps of the proof. It now follows from (24) and the monotonicity of f that

$$\bar{P}_{e,u}(\mathbf{y}, \mathbf{z}) \leq f\left(e^{n\epsilon_1(n)} \cdot \exp_2[u(\mathbf{y}, \mathbf{z})]\right). \tag{25}$$

For later use, we also have

$$\bar{P}_{e,u}(\mathbf{z}) \triangleq \sum_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{z}) \bar{P}_{e,u}(\mathbf{y}, \mathbf{z}) \tag{26}$$

$$\leq \sum_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{z}) f\left(e^{n\epsilon_1(n)} \cdot \exp_2[u(\mathbf{y}, \mathbf{z})]\right). \tag{27}$$

We next move on to derive a matching lower bound to $\bar{P}_{e,o}$. Similarly, as in [11], we will need to refer to an auxiliary threshold decoder (in the terminology of [11]), which is a slightly more conservative version of the ML decoder. Specifically, for a given threshold parameter, $\alpha > 1$, this decoder outputs the message m with the property that $P(\mathbf{z}|\mathbf{y}_m) > \alpha \cdot P(\mathbf{z}|\mathbf{y}_{m'})$ for all $m' \neq m$, and declares an error if no such m exists. Accordingly, let $\bar{P}_{e,t}(\mathbf{y}, \mathbf{z})$ denote the conditional average error probability of the threshold decoder, given (\mathbf{y}, \mathbf{z}) , i.e.,

$$\bar{P}_{e,t}(\mathbf{y}, \mathbf{z}) = f(P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})]), \tag{28}$$

where

$$\mathcal{E}_t(\mathbf{y}, \mathbf{z}) = \{\mathbf{y}' : P(\mathbf{z}|\mathbf{y}') \geq \alpha^{-1} P(\mathbf{z}|\mathbf{y})\}. \tag{29}$$

As in Lemma 2 of [11], here too, the next lemma (proved in the appendix) asserts that the performance of the threshold decoder cannot be much worse than that of the ML decoder, provided that α is not too large. In particular, if $\alpha = \alpha_n$ grows subexponentially with n , then the threshold decoder has the same error exponent as that of the ML decoder.

Lemma 1 *Define*

$$\bar{P}_{e,t}(\mathbf{z}) = \sum_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{z}) f(P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})]) \tag{30}$$

$$\bar{P}_{e,o}(\mathbf{z}) = \sum_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{z}) f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]). \quad (31)$$

Then, under the positivity assumption (7),

$$\bar{P}_{e,t}(\mathbf{z}) \leq \left\{ \alpha \left[n \ln \left(\frac{1}{\pi_{\min} \cdot |\Theta| \cdot |\Omega|} \right) + 1 \right] + 1 \right\} \cdot \bar{P}_{e,o}(\mathbf{z}) \quad (32)$$

for every $\mathbf{z} \in \mathcal{Z}^n$.

It should be noted that assumption (7) is essentially not needed for the above Lemma. What is really needed is that the smallest $P(\mathbf{y})$, across all $\mathbf{y} \in \mathcal{Y}^n$ with $P(\mathbf{y}) > 0$, would not decay faster than exponentially with n . But owing to (6), one can easily see that $P(\mathbf{y}) \geq \pi_+^n$, where π_+ is the smallest positive $\pi(y, \theta, \omega|\theta', \omega')$. We are using (7) nonetheless, because we make this assumption anyway (as it is needed elsewhere), and then the upper bound given by the lemma is slightly tighter.

On the basis of Lemma 1, any lower bound on $\bar{P}_{e,t}$ in terms of $\bar{P}_{e,u}$, would immediately yield a lower bound $\bar{P}_{e,o}$ in terms of $\bar{P}_{e,u}$, as desired. Accordingly, the next step would be to lower bound $\bar{P}_{e,t}$. This in turn will be done by lower bounding $P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})]$ (for a certain choice of the threshold α , to be defined) in terms of $P[\mathcal{E}_1(\mathbf{y}, \mathbf{z})]$, for a certain $\mathcal{E}_1(\mathbf{y}, \mathbf{z}) \subseteq \mathcal{E}_t(\mathbf{y}, \mathbf{z})$ to be specified shortly.

First observe that, similarly as in eq. (6),

$$P(\mathbf{y}, \mathbf{z}) = \sum_{\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\omega}, \mathbf{x}} \prod_{i=1}^n [G(x_i, \omega_i | \omega_{i-1}) V(y_i, \theta_i | x_i, \theta_{i-1}) W(z_i, \sigma_i | x_i, \sigma_{i-1})] \quad (33)$$

$$= \sum_{\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\omega}} \prod_{i=1}^n \sum_x [G(x, \omega_i | \omega_{i-1}) V(y_i, \theta_i | x, \theta_{i-1}) W(z_i, \sigma_i | x, \sigma_{i-1})] \quad (34)$$

$$= \sum_{\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\omega}} \prod_{i=1}^n \Pi(y_i, z_i, \theta_i, \sigma_i, \omega_i | \theta_{i-1}, \sigma_{i-1}, \omega_{i-1}) \quad (35)$$

where we have defined $\Pi(y, z, \theta, \sigma, \omega | \theta', \sigma', \omega') = \sum_x G(x, \omega | \omega') V(y, \theta | x, \theta') W(z, \sigma | x, \sigma')$. We will henceforth use the following notation for two positive integers i and j , where $j > i$:

$$\begin{aligned} & \Pi(y_i^j, z_i^j, \theta_j, \sigma_j, \omega_j | \theta_{i-1}, \sigma_{i-1}, \omega_{i-1}) \\ &= \sum_{\theta_i^{j-1}} \sum_{\sigma_i^{j-1}} \sum_{\omega_i^{j-1}} \prod_{k=i}^j \Pi(y_k, z_k, \theta_k, \sigma_k, \omega_k | \theta_{k-1}, \sigma_{k-1}, \omega_{k-1}) \end{aligned} \quad (36)$$

and

$$\pi(y_i^j, \theta_j, \omega_j | \theta_{i-1}, \omega_{i-1}) = \sum_{\theta_i^{j-1}} \sum_{\omega_i^{j-1}} \prod_{k=i}^j \pi(y_k, \theta_k, \omega_k | \theta_{k-1}, \omega_{k-1}). \quad (37)$$

Next, define

$$\mathbf{t} \triangleq \{(\theta_i, \sigma_i, \omega_i) : i = n_0, n_1, \dots, n_{c-1}\}, \quad (38)$$

$$\mathbf{s} \triangleq \{(\theta_i, \omega_i) : i = n_0, n_1, \dots, n_{c-1}\}, \quad (39)$$

where $\{n_i\}$ are phrase boundaries, as defined at the beginning of Section 3, for a given (\mathbf{y}, \mathbf{z}) . Now, for the same (\mathbf{y}, \mathbf{z}) , let

$$\hat{\mathbf{t}} = \arg \max_{\mathbf{t}} P(\mathbf{y}, \mathbf{z}, \mathbf{t}) = \arg \max_{\mathbf{t}} \prod_{i=0}^{c-1} \Pi(y_{n_i+1}^{n_{i+1}}, z_{n_i+1}^{n_{i+1}}, \theta_{n_i+1}, \sigma_{n_i+1}, \omega_{n_i+1} | \theta_{n_i}, \sigma_{n_i}, \omega_{n_i}) \quad (40)$$

$$\tilde{\mathbf{s}} = \arg \max_{\mathbf{s}} P(\mathbf{y}, \mathbf{s}) = \arg \max_{\mathbf{s}} \prod_{i=0}^{c-1} \pi(y_{n_i+1}^{n_{i+1}}, \theta_{n_i+1}, \omega_{n_i+1} | \theta_{n_i}, \omega_{n_i}). \quad (41)$$

We denote the components of $\hat{\mathbf{t}}$ and $\tilde{\mathbf{s}}$ by $\{(\hat{\theta}_{n_i}, \hat{\sigma}_{n_i}, \hat{\omega}_{n_i})\}$ and $\{(\tilde{\theta}_{n_i}, \tilde{\omega}_{n_i})\}$, respectively. Denoting $K = |\Theta \times \Sigma \times \Omega|$, it is obvious that $P(\mathbf{y}, \mathbf{z}, \hat{\mathbf{t}}) \geq K^{-c} P(\mathbf{y}, \mathbf{z})$, and a similar relation holds between $P(\mathbf{y}, \tilde{\mathbf{s}})$ and $P(\mathbf{y})$. For the given pair (\mathbf{y}, \mathbf{z}) , let

$$\mathcal{E}_1(\mathbf{y}, \mathbf{z}) \triangleq \left\{ \mathbf{y}' : P(\mathbf{y}', \mathbf{z}, \hat{\mathbf{t}}) = P(\mathbf{y}, \mathbf{z}, \hat{\mathbf{t}}), P(\mathbf{y}', \tilde{\mathbf{s}}) = P(\mathbf{y}, \tilde{\mathbf{s}}) \right\}. \quad (42)$$

Owing to assumption (7), it is shown in the appendix (similarly as in [22, eq. (A.7)]) that

$$P(\mathbf{y}') \leq P(\mathbf{y}', \tilde{\mathbf{s}}) \cdot \left(\frac{|\Theta \times \Omega|}{\pi_{\min}^2} \right)^c \leq P(\mathbf{y}', \tilde{\mathbf{s}}) \cdot \left(\frac{K}{\pi_{\min}^2} \right)^c, \quad (43)$$

and so, for $\mathbf{y}' \in \mathcal{E}_1(\mathbf{y}, \mathbf{z})$, the chain of inequalities,

$$\left(\frac{K}{\pi_{\min}^2} \right)^c \cdot P(\mathbf{z} | \mathbf{y}') = \left(\frac{K}{\pi_{\min}^2} \right)^c \cdot \frac{P(\mathbf{y}', \mathbf{z})}{P(\mathbf{y}')} \quad (44)$$

$$\geq \left(\frac{K}{\pi_{\min}^2} \right)^c \frac{P(\mathbf{y}', \mathbf{z}, \hat{\mathbf{t}})}{(K/\pi_{\min}^2)^c P(\mathbf{y}', \tilde{\mathbf{s}})} \quad (45)$$

$$= \frac{P(\mathbf{y}', \mathbf{z}, \hat{\mathbf{t}})}{P(\mathbf{y}', \tilde{\mathbf{s}})} \quad (46)$$

$$= \frac{P(\mathbf{y}, \mathbf{z}, \hat{\mathbf{t}})}{P(\mathbf{y}, \tilde{\mathbf{s}})} \quad (47)$$

$$\geq K^{-c} \frac{P(\mathbf{y}, \mathbf{z})}{P(\mathbf{y})} \quad (48)$$

$$= K^{-c} P(\mathbf{z} | \mathbf{y}), \quad (49)$$

implies that

$$\mathcal{E}_1(\mathbf{y}, \mathbf{z}) \subseteq \{ \mathbf{y}' : P(\mathbf{z} | \mathbf{y}') \geq (K/\pi_{\min})^{-2c} P(\mathbf{z} | \mathbf{y}) \} \quad (50)$$

$$\subseteq \{ \mathbf{y}' : P(\mathbf{z} | \mathbf{y}') \geq (K/\pi_{\min})^{-2\bar{c}_n} P(\mathbf{z} | \mathbf{y}) \} \quad (51)$$

$$= \mathcal{E}_t(\mathbf{y}, \mathbf{z}) \quad \text{with the choice } \alpha = (K/\pi_{\min})^{2\bar{c}_n} \quad (52)$$

where

$$\bar{c}_n \triangleq \frac{n \log |\mathcal{Y} \times \mathcal{Z}|}{(1 - \varepsilon_n) \log n}, \quad (53)$$

with $\varepsilon_n \rightarrow 0$ as $n \rightarrow 0$, so that \bar{c}_n serves as a uniform upper bound to $c \equiv c(\mathbf{y}, \mathbf{z})$ for every $(\mathbf{y}, \mathbf{z}) \in \mathcal{Y} \times \mathcal{Z}^n$, according to [21, eq. (6)]. Thus,

$$P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})] = \sum_{\mathbf{y}' \in \mathcal{E}_t(\mathbf{y}, \mathbf{z})} P(\mathbf{y}') \quad (54)$$

$$\geq \sum_{\mathbf{y}' \in \mathcal{E}_1(\mathbf{y}, \mathbf{z})} P(\mathbf{y}') \quad (55)$$

$$\geq \sum_{\mathbf{y}' \in \mathcal{E}_1(\mathbf{y}, \mathbf{z})} P(\mathbf{y}', \tilde{\mathbf{s}}) \quad (56)$$

$$= \sum_{\mathbf{y}' \in \mathcal{E}_1(\mathbf{y}, \mathbf{z})} P(\mathbf{y}, \tilde{\mathbf{s}}) \quad (57)$$

$$= |\mathcal{E}_1(\mathbf{y}, \mathbf{z})| \cdot P(\mathbf{y}, \tilde{\mathbf{s}}) \quad (58)$$

$$\geq K^{-c} \cdot |E_1(\mathbf{y}, \mathbf{z})| \cdot P(\mathbf{y}) \quad (59)$$

$$\geq K^{-\bar{c}_n} \cdot |E_1(\mathbf{y}, \mathbf{z})| \cdot P(\mathbf{y}). \quad (60)$$

Now, let $\mathcal{T}(\mathbf{y}|\mathbf{z}, \hat{\mathbf{t}}, \tilde{\mathbf{s}})$ denote the set of all $\mathbf{y}' \in \mathcal{Y}^n$ that are obtained from \mathbf{y} by permuting \mathbf{y} -phrases, $\{y_{n_i+1}^{n_{i+1}}\}$, that are: (i) aligned to the same \mathbf{z} -phrases, $z_{n_i+1}^{n_{i+1}}$, (ii) of the same length, (iii) begin at the same states, of both $\hat{t}_i = (\hat{\theta}_{n_i}, \hat{\sigma}_{n_i}, \hat{\omega}_{n_i})$ and $\tilde{s}_i = (\tilde{\theta}_{n_i}, \tilde{\omega}_{n_i})$, and (iv) end at the same states of both $\hat{t}_{i+1} = (\hat{\theta}_{n_{i+1}}, \hat{\sigma}_{n_{i+1}}, \hat{\omega}_{n_{i+1}})$ and $\tilde{s}_{i+1} = (\tilde{\theta}_{n_{i+1}}, \tilde{\omega}_{n_{i+1}})$. Clearly, $\mathcal{T}(\mathbf{y}|\mathbf{z}, \hat{\mathbf{t}}, \tilde{\mathbf{s}}) \subseteq \mathcal{E}_1(\mathbf{y}, \mathbf{z})$, and so, $P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})]$ is further lower bounded by

$$P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})] \geq K^{-\bar{c}_n} |\mathcal{T}(\mathbf{y}|\mathbf{z}, \hat{\mathbf{t}}, \tilde{\mathbf{s}})| \cdot P(\mathbf{y}). \quad (61)$$

Now, according to Lemma 1 of [20],

$$|\mathcal{T}(\mathbf{y}|\mathbf{z}, \hat{\mathbf{t}}, \tilde{\mathbf{s}})| \geq \exp_2\{v(\mathbf{y}, \mathbf{z}) - n\epsilon'_2(n)\}, \quad (62)$$

where

$$\epsilon'_2(n) = \frac{\bar{c}_n}{n} \cdot \log(|\Theta|^4 \cdot |\Omega|^4 \cdot |\Sigma|^2 e) \quad (63)$$

$$= \frac{\log(|\mathcal{Y}| \cdot |\mathcal{Z}|)}{(1 - \varepsilon_n) \log n} \cdot \log(|\Theta|^4 \cdot |\Omega|^4 \cdot |\Sigma|^2 e) \quad (64)$$

$$= O\left(\frac{1}{\log n}\right). \quad (65)$$

Thus,

$$P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})] \geq K^{-\bar{c}_n} P(\mathbf{y}) \cdot \exp_2\{v(\mathbf{y}, \mathbf{z}) - n\epsilon'_2(n)\} \stackrel{\Delta}{=} \exp_2\{u(\mathbf{y}, \mathbf{z}) - n\epsilon_2(n)\} \quad (66)$$

where

$$\epsilon_2(n) = \epsilon'_2(n) + \frac{\bar{c}_n \log K}{n} \quad (67)$$

$$\leq \epsilon'_2(n) + \frac{\log(|\mathcal{Y}| \cdot |\mathcal{Z}|) \cdot \log K}{(1 - \varepsilon_n) \log n} \quad (68)$$

$$= O\left(\frac{1}{\log n}\right), \quad (69)$$

and so,

$$P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})] \geq \exp_2\{u(\mathbf{y}, \mathbf{z}) - n\epsilon_2(n)\}. \quad (70)$$

To complete the proof, we use the first part of Lemma 1 of [11], which asserts that for every $a, b \in [0, 1]$, $f(a)/f(b) \leq \max\{1, a/b\}$, and so,

$$\frac{f(P[\mathcal{E}_u(\mathbf{y}, \mathbf{z})])}{f(P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})])} \leq \max\left\{1, \frac{P[\mathcal{E}_u(\mathbf{y}, \mathbf{z})]}{P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})]}\right\} \quad (71)$$

$$\leq \max\left\{1, \frac{\exp_2\{u(\mathbf{y}, \mathbf{z}) + n\epsilon_1(n)\}}{\exp_2\{u(\mathbf{y}, \mathbf{z}) - n\epsilon_2(n)\}}\right\} \quad (72)$$

$$\leq e^{n[\epsilon_1(n) + \epsilon_2(n)]}, \quad (73)$$

where in the second inequality, we have used eqs. (24) and (70). Now, referring to Lemma 1, let us define

$$\epsilon_3(n) = \frac{1}{n} \log \left\{ \left(\frac{K}{\pi_{\min}} \right)^{2\bar{c}_n} \left[n \ln \left(\frac{1}{\pi_{\min} |\Theta \times \Sigma|} \right) + 1 \right] + 1 \right\} \quad (74)$$

$$= O\left(\frac{1}{\log n}\right). \quad (75)$$

Then,

$$\bar{P}_{e,o}(\mathbf{z}) \geq e^{-n\epsilon_3(n)} \bar{P}_{e,t}(\mathbf{z}) \quad (\text{by Lemma 1}) \quad (76)$$

$$= e^{-n\epsilon_3(n)} \sum_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{z}) f(P[\mathcal{E}_t(\mathbf{y}, \mathbf{z})]) \quad (77)$$

$$\geq e^{-n[\epsilon_1(n) + \epsilon_2(n) + \epsilon_3(n)]} \sum_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{z}) f(P[\mathcal{E}_u(\mathbf{y}, \mathbf{z})]) \quad (78)$$

$$= e^{-n[\epsilon_1(n) + \epsilon_2(n) + \epsilon_3(n)]} \bar{P}_{e,u}(\mathbf{z}). \quad (79)$$

Finally, upon averaging both sides over $\{\mathbf{z}\}$, we complete the proof of Theorem 1, with

$$\epsilon(n) \triangleq \epsilon_1(n) + \epsilon_2(n) + \epsilon_3(n), \quad (80)$$

which is $O((\log \log n)/\log n)$ since $\epsilon_1(n)$ is such.

Appendix

A1. Proof of Lemma 1

Let us define

$$\Delta(\mathbf{y}, \mathbf{z}) \triangleq \{\mathbf{y}' : M_o(\mathbf{y}', \mathbf{z}) > M_o(\mathbf{y}, \mathbf{z}), P(\mathbf{z}|\mathbf{y}') \geq \alpha^{-1} P(\mathbf{z}|\mathbf{y})\} \quad (\text{A.1})$$

$$= \{\mathbf{y}' : M_o(\mathbf{y}', \mathbf{z}) > M_o(\mathbf{y}, \mathbf{z}), P(\mathbf{y})P(\mathbf{y}'|\mathbf{z}) \geq \alpha^{-1}P(\mathbf{y}')P(\mathbf{y}|\mathbf{z})\}, \quad (\text{A.2})$$

so that $\mathcal{E}_t(\mathbf{y}, \mathbf{z})$ is given by the disjoint union of $\mathcal{E}_o(\mathbf{y}, \mathbf{z})$ and $\Delta(\mathbf{y}, \mathbf{z})$. Then the average conditional error probabilities given \mathbf{z} are

$$\bar{P}_{e,o}(\mathbf{z}) = \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{z}) f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]) \quad (\text{A.3})$$

$$\bar{P}_{e,t}(\mathbf{z}) = \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{z}) f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})] + P[\Delta(\mathbf{y}, \mathbf{z})]) \quad (\text{A.4})$$

$$\leq \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{z}) \left(\frac{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})] + P[\Delta(\mathbf{y}, \mathbf{z})]}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]} \right) f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]), \quad (\text{A.5})$$

where in the last step, we have used the first part of Lemma 1 from [11] (see also [7]). Now, let us define

$$r(\mathbf{y}, \mathbf{z}) \triangleq \sum_{\mathbf{y}' \in \mathcal{E}_o(\mathbf{y}, \mathbf{z})} P(\mathbf{y}'|\mathbf{z}). \quad (\text{A.6})$$

Then,

$$P(\mathbf{y}) = \sum_{\mathbf{y}'} P(\mathbf{y})P(\mathbf{y}'|\mathbf{z}) \quad (\text{A.7})$$

$$\geq \sum_{\mathbf{y}' \in \mathcal{E}_o(\mathbf{y}, \mathbf{z})} P(\mathbf{y})P(\mathbf{y}'|\mathbf{z}) + \sum_{\mathbf{y}' \in \Delta(\mathbf{y}, \mathbf{z})} P(\mathbf{y})P(\mathbf{y}'|\mathbf{z}) \quad (\text{A.8})$$

$$= P(\mathbf{y})r(\mathbf{y}, \mathbf{z}) + \sum_{\mathbf{y}' \in \Delta(\mathbf{y}, \mathbf{z})} P(\mathbf{y})P(\mathbf{y}'|\mathbf{z}) \quad (\text{A.9})$$

$$\geq P(\mathbf{y})r(\mathbf{y}, \mathbf{z}) + \frac{1}{\alpha} \sum_{\mathbf{y}' \in \Delta(\mathbf{y}, \mathbf{z})} P(\mathbf{y}')P(\mathbf{y}|\mathbf{z}) \quad (\text{A.10})$$

$$= P(\mathbf{y})r(\mathbf{y}, \mathbf{z}) + \frac{P(\mathbf{y}|\mathbf{z})}{\alpha} P[\Delta(\mathbf{y}, \mathbf{z})], \quad (\text{A.11})$$

and so,

$$P(\mathbf{y}|\mathbf{z})P[\Delta(\mathbf{y}, \mathbf{z})] \leq \alpha P(\mathbf{y})[1 - r(\mathbf{y}, \mathbf{z})]. \quad (\text{A.12})$$

We then have

$$\bar{P}_{e,t}(\mathbf{z}) - \bar{P}_{e,o}(\mathbf{z}) \quad (\text{A.13})$$

$$\leq \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{z}) \frac{P[\Delta(\mathbf{y}, \mathbf{z})]}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]} f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]) \quad (\text{A.14})$$

$$\leq \alpha \cdot \sum_{\mathbf{y}} \frac{P(\mathbf{y})[1 - r(\mathbf{y}, \mathbf{z})]}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]} f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]) \quad (\text{A.15})$$

$$= \alpha \cdot \sum_{\mathbf{y}} \sum_{\{\mathbf{y}': M_o(\mathbf{y}', \mathbf{z}) > M_o(\mathbf{y}, \mathbf{z})\}} \frac{P(\mathbf{y})P(\mathbf{y}'|\mathbf{z})}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]} f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]) \quad (\text{A.16})$$

$$\stackrel{(a)}{=} \alpha \cdot \sum_{\mathbf{y}'} \sum_{\{\mathbf{y}: M_o(\mathbf{y}', \mathbf{z}) > M_o(\mathbf{y}, \mathbf{z})\}} \frac{P(\mathbf{y})P(\mathbf{y}'|\mathbf{z})}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]} f(P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]) \quad (\text{A.17})$$

$$\stackrel{(b)}{\leq} \alpha \cdot \sum_{\mathbf{y}'} \sum_{\{\mathbf{y}: M_o(\mathbf{y}', \mathbf{z}) > M_o(\mathbf{y}, \mathbf{z})\}} \frac{P(\mathbf{y})P(\mathbf{y}'|\mathbf{z})}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]} f(P[\mathcal{E}_o(\mathbf{y}', \mathbf{z})]) \quad (\text{A.18})$$

$$\leq \alpha \cdot \sum_{\mathbf{y}'} P(\mathbf{y}'|\mathbf{z}) f(P[\mathcal{E}_o(\mathbf{y}', \mathbf{z})]) \cdot \sum_{\mathbf{y}} \frac{P(\mathbf{y})}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]} \quad (\text{A.19})$$

$$= \alpha \cdot \bar{P}_{e,o}(\mathbf{z}) \cdot \sum_{\mathbf{y} \in \mathcal{Y}^n} \frac{P(\mathbf{y})}{P[\mathcal{E}_o(\mathbf{y}, \mathbf{z})]}, \quad (\text{A.20})$$

where in (a) we have interchanged the order of the summation and in (b), we have used the monotonicity of f together with the fact that $\mathcal{E}_o(\mathbf{y}, \mathbf{z}) \subseteq \mathcal{E}_o(\mathbf{y}', \mathbf{z})$ whenever $M_o(\mathbf{y}', \mathbf{z}) > M_o(\mathbf{y}, \mathbf{z})$. To complete the proof, it remains to show then that for any \mathbf{z} ,

$$L_n(\mathbf{z}) \triangleq \sum_{\mathbf{y} \in \mathcal{Y}^n} \frac{P(\mathbf{y})}{P[\mathcal{E}(\mathbf{y}, \mathbf{z})]} = \sum_{\mathbf{y} \in \mathcal{Y}^n} \frac{P(\mathbf{y})}{\sum_{\{\mathbf{y}': M_o(\mathbf{y}', \mathbf{z}) \leq M_o(\mathbf{y}, \mathbf{z})\}} P(\mathbf{y}')} \quad (\text{A.21})$$

cannot exceed $1 + n \ln[1/(\pi_{\min}|\Theta \times \Omega|)]$. For the given \mathbf{z} , consider the ordering of all members of \mathcal{Y}^n (not only those in \mathcal{C}) according to the ranking function $M_o(\mathbf{y}, \mathbf{z})$, i.e.,

$$P(\mathbf{z}|\mathbf{y}[1]) \geq P(\mathbf{z}|\mathbf{y}[2]) \geq \dots \geq P(\mathbf{z}|\mathbf{y}[N]), \quad N = |\mathcal{Y}|^n \quad (\text{A.22})$$

and let us denote $a_i = P(\mathbf{y}[i])$, $A_i = \sum_{j=1}^i a_j$, $i = 1, \dots, N$. Then, using the facts that $A_1 = a_1 = P(\mathbf{y}[1])$ and $A_N = 1$, as well as the inequality

$$\ln(1+u) \equiv -\ln\left(1 - \frac{u}{1+u}\right) \geq \frac{u}{1+u}, \quad (\text{A.23})$$

we have

$$L_n(\mathbf{z}) = \sum_{i=1}^N \frac{a_i}{A_i} \quad (\text{A.24})$$

$$= 1 + \sum_{i=2}^N \frac{a_i}{A_{i-1} + a_i} \quad (\text{A.25})$$

$$= 1 + \sum_{i=2}^N \frac{a_i/A_{i-1}}{1 + a_i/A_{i-1}} \quad (\text{A.26})$$

$$\leq 1 + \sum_{i=2}^N \ln\left(1 + \frac{a_i}{A_{i-1}}\right) \quad (\text{A.27})$$

$$= 1 + \sum_{i=2}^N \ln\left(\frac{A_{i-1} + a_i}{A_{i-1}}\right) \quad (\text{A.28})$$

$$= 1 + \sum_{i=2}^N \ln\left(\frac{A_i}{A_{i-1}}\right) \quad (\text{A.29})$$

$$= 1 + \ln\left(\frac{A_N}{A_1}\right) \quad (\text{A.30})$$

$$= \ln \left[\frac{1}{P(\mathbf{y}[1])} \right] + 1 \quad (\text{A.31})$$

$$\leq \ln \left[\frac{1}{(\pi_{\min} \cdot |\Theta| \cdot |\Omega|)^n} \right] + 1 \quad (\text{A.32})$$

$$= n \ln \left(\frac{1}{\pi_{\min} \cdot |\Theta| \cdot |\Omega|} \right) + 1, \quad (\text{A.33})$$

where we have used the assumption (7), which implies that $P(\mathbf{y}) \geq (\pi_{\min} \cdot |\Theta| \cdot |\Omega|)^n$ for all \mathbf{y} . This completes the proof of Lemma 1.

A.2 Proof of Eq. (43)

We next show that for every \mathbf{y} and \mathbf{s} ,

$$P(\mathbf{y}) \leq P(\mathbf{y}, \mathbf{s}) \cdot \left(\frac{|\Theta \times \Omega|}{\pi_{\min}^2} \right)^c. \quad (\text{A.34})$$

For the sake of brevity, let us denote $\zeta_i = (\theta_i, \omega_i)$ (so that $s_i = \zeta_{n_i}$). Now,

$$P(\mathbf{y}, \mathbf{s}) = \prod_{i=0}^{c-1} \pi(y_{n_i+1}^{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i}). \quad (\text{A.35})$$

But

$$\pi(y_{n_i+1}^{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i}) = \sum_{\zeta_{n_i+1}^{n_i+1-1}} \prod_{t=n_i+1}^{n_i+1} \pi(y_t, \zeta_t | \zeta_{t-1}) \quad (\text{A.36})$$

$$\begin{aligned} &= \sum_{\zeta_{n_i+1}} \pi(y_{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i}) \times \\ &\quad \sum_{\zeta_{n_i+2}^{n_i+1-2}} \prod_{t=n_i+2}^{n_i+1-1} \pi(y_t, \zeta_t | \zeta_{t-1}) \times \\ &\quad \sum_{\zeta_{n_i+1-1}} \pi(y_{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i+1-1}) \end{aligned} \quad (\text{A.37})$$

$$\geq \pi_{\min}^2 \sum_{\zeta_{n_i+1}^{n_i+1-1}} \prod_{t=n_i+2}^{n_i+1-1} \pi(y_t, \zeta_t | \zeta_{t-1}), \quad (\text{A.38})$$

where we have assumed that $n_i + 2 \leq n_{i+1} - 1$, which means that the phrase length must be at least three,⁹ and where we have lower bounded both $\pi(y_{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i})$ and $\pi(y_{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i+1-1})$ by π_{\min} . Similarly, since both $\pi(y_{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i})$ and $\pi(y_{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i+1-1})$ are upper bounded by unity, we have

$$\pi(y_{n_i+1}^{n_i+1}, \zeta_{n_i+1} | \zeta_{n_i}) \leq \sum_{\zeta_{n_i+1}^{n_i+1-1}} \prod_{t=n_i+2}^{n_i+1-1} \pi(y_t, \zeta_t | \zeta_{t-1}). \quad (\text{A.39})$$

⁹This assumption does not affect the generality, as the number of phrases of length shorter than three cannot exceed $|\mathcal{Y} \times \mathcal{Z}| + |\mathcal{Y} \times \mathcal{Z}|^2$, which is fixed and hence negligible compared to the total number of phrases for large n .

Since the expression

$$\sum_{\zeta_{n_i+1}}^{n_{i+1}-1} \prod_{t=n_i+2}^{n_{i+1}-1} \pi(y_t, \zeta_t | \zeta_{t-1})$$

depends neither on ζ_{n_i} nor on $\zeta_{n_{i+1}}$, it follows that for any ζ_{n_i} , ζ'_{n_i} , $\zeta_{n_{i+1}}$, and $\zeta'_{n_{i+1}}$,

$$\pi_{\min}^2 \leq \frac{\pi(y_{n_i+1}^{n_{i+1}}, \zeta'_{n_{i+1}} | \zeta'_{n_i})}{\pi(y_{n_i+1}^{n_{i+1}}, \zeta_{n_{i+1}} | \zeta_{n_i})} \leq \frac{1}{\pi_{\min}^2}, \quad (\text{A.40})$$

and so,

$$P(\mathbf{y}) = \sum_{\mathbf{s}'} P(\mathbf{y}, \mathbf{s}') \quad (\text{A.41})$$

$$= P(\mathbf{y}, \mathbf{s}) \sum_{\mathbf{s}'} \frac{P(\mathbf{y}, \mathbf{s}')}{P(\mathbf{y}, \mathbf{s})} \quad (\text{A.42})$$

$$= P(\mathbf{y}, \mathbf{s}) \sum_{\mathbf{s}'} \prod_{i=0}^{c-1} \frac{\pi(y_{n_i+1}^{n_{i+1}}, \zeta'_{n_{i+1}} | \zeta'_{n_i})}{\pi(y_{n_i+1}^{n_{i+1}}, \zeta_{n_{i+1}} | \zeta_{n_i})} \quad (\text{A.43})$$

$$\leq P(\mathbf{y}, \mathbf{s}) \sum_{\mathbf{s}'} \prod_{i=0}^{c-1} \frac{1}{\pi_{\min}^2} \quad (\text{A.44})$$

$$= P(\mathbf{y}, \mathbf{s}) \cdot \left(\frac{|\Omega \times \Theta|}{\pi_{\min}^2} \right)^c, \quad (\text{A.45})$$

which completes the proof of eq. (43).

References

- [1] L. E. Baum and T. Petrie, "Statistical inference for probabilistic functions of finite state Markov chains," *Ann. Math. Statist.*, vol. 37, pp. 1554–1563, 1966.
- [2] I. Csiszár, "Linear codes for sources and source networks: error exponents, universal coding," *IEEE Trans. Inform. Theory*, vol. IT-28, no. 4, pp. 585–592, July 1982.
- [3] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, Cambridge University Press, 2011.
- [4] G. Dasarathy and S. C. Draper, "On reliability of content identification from databases based on noisy queries," *The 2011 IEEE Proc. International Symposium on Information Theory (ISIT 2011)*, pp. 1066–1070, St. Petersburg, Russia, July–August 2011.
- [5] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Roy. Statist. Soc.*, vol. B39, pp. 1–39, 1977.
- [6] Y. Ephraim and N. Merhav, "Hidden Markov processes," *IEEE Trans. Inform. Theory*, special issue in memory of Aaron D. Wyner, June 2002.
- [7] M. Feder and A. Lapidoth, "Universal decoding for channels with memory," *IEEE Trans. Inform. Theory*, vol. 44, no. 5, pp. 1726–1745, September 1998.
- [8] M. Feder and N. Merhav, "Universal composite hypothesis testing: a competitive minimax approach," *IEEE Trans. Inform. Theory*, special issue in memory of Aaron D. Wyner, vol. 48, no. 6, pp. 1504–1517, June 2002.
- [9] V. D. Goppa, "Nonprobabilistic mutual information without memory," *Probl. Cont. Information Theory*, vol. 4, pp. 97–102, 1975.
- [10] T. Ignatenko and F. M. J. Willems, "Biometric security from an information-theoretical perspective," *Foundations and Trends in Communications and Information Theory*, vol. 7, nos. 2–3, pp. 135–316.
- [11] A. Lapidoth and J. Ziv, "On the universality of the LZ-based noisy channels decoding algorithm," *IEEE Trans. Inform. Theory*, vol. 44, no. 5, pp. 1746–1755, September 1998.
- [12] Y. Lomnitz and M. Feder, "Communication over individual channels – a general framework," *IEEE Trans. Inform. Theory*, vol. 57, no. 11, pp. 7333–7358, November 2011.
- [13] N. Merhav, "Universal decoding for memoryless Gaussian channels with a deterministic interference," *IEEE Trans. Inform. Theory*, vol. 39, no. 4, pp. 1261–1269, July 1993.
- [14] N. Merhav, "Universal decoding for arbitrary channels relative to a given family of decoding metrics," *IEEE Trans. Inform. Theory*, vol. 59, no. 9, pp. 5566–5576, September 2013.
- [15] T. Petrie, "Probabilistic functions of finite state Markov chains," *Ann. Math. Statist.*, vol. 40, no. 1, pp. 97–115, 1969.

- [16] L. R. Rabiner, "A tutorial of hidden Markov models and selected applications in speech recognition", *Proceedings of the IEEE*, vol. 77, no. 2, February 1989.
- [17] E. Tuncel, "Capacity/storage tradeoff in high-dimensional identification systems," *IEEE Trans. Inform. Theory*, vol. 55, no. 5, pp. 2097–2106, May 2009.
- [18] F. Willems, T. Kalker, J. Goseling, and J.-P. Linnartz, "On the capacity of a biometrical identification system," *The 2003 IEEE Proc. International Symposium on Information Theory (ISIT 2003)*, p. 82, Yokohama, Japan, June–July 2003.
- [19] F. Willems, T. Kalker, S. Baggen, and J.-P. Linnartz, "On the capacity of a biometrical identification system," (unknown year) available on-line at:
<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.74.9512&rep=rep1&type=pdf>
- [20] J. Ziv, "Universal decoding for finite-state channels," *IEEE Trans. Inform. Theory*, vol. IT-31, no. 4, pp. 453–460, July 1985.
- [21] J. Ziv and A. Lempel, "Compression of individual sequences via variable-rate coding," *IEEE Trans. Inform. Theory*, vol. IT-24, no. 5, pp. 530–536, September 1978.
- [22] J. Ziv and N. Merhav, "Estimating the number of states of a finite-state source," *IEEE Trans. Inform. Theory*, vol. 38, no. 1, pp. 61–65, January 1992.